

20 ABSTRACT

21 Contact tracing requires reliable identification of closely related bacterial isolates. When we noticed
22 the reporting of artefactual variation between *M. tuberculosis* isolates during routine next
23 generation sequencing of *Mycobacterium spp*, we investigated its basis in 2,018 consecutive *M.*
24 *tuberculosis* isolates. In the routine process used, clinical samples were decontaminated and
25 inoculated into broth cultures; from positive broth cultures DNA was extracted, sequenced, reads
26 mapped, and consensus sequences determined. We investigated the process of consensus
27 sequence determination, which selects the most common nucleotide at each position. Having
28 determined the high-quality read depth and depth of minor variants across 8,006 *M. tuberculosis*
29 genomic regions, we quantified the relationship between the minor variant depth and the amount of
30 non-*Mycobacterium* bacterial DNA, which originates from commensal microbes killed during sample
31 decontamination. In the presence of non-*Mycobacterium* bacterial DNA, we found significant
32 increases in minor variant frequencies of more than 1.5 fold in 242 regions covering 5.1% of the *M.*
33 *tuberculosis* genome. Included within these were four high variation regions strongly influenced by
34 the amount of non-*Mycobacterium* bacterial DNA. Excluding these four regions from pairwise
35 distance comparisons reduced biologically implausible variation from 5.2% to 0% in an independent
36 validation set derived from 226 individuals. Thus, we have demonstrated an approach identifying
37 critical genomic regions contributing to clinically relevant artefactual variation in bacterial similarity
38 searches. The approach described monitors the outputs of the complex multi-step laboratory and
39 bioinformatics process, allows periodic process adjustments, and will have application to quality
40 control of routine bacterial genomics.

41

42

43 INTRODUCTION

44 Identifying closely related bacterial isolates is required for clinical and epidemiological purposes (1-
45 3). Most published approaches using short-read next generation sequencing (NGS) rely on mapping
46 to a high-quality reference sequence followed by consensus base calling (1-8). A known problem
47 with this approach concerns the existence in many bacterial genomes of 'hard-to-map' regions
48 which are either repeated within the genome, or which contain regions of low sequence complexity.
49 High-confidence mapping of short reads to such regions is difficult or impossible, and so determining
50 the consensus sequence of these regions is difficult. One approach to managing this problem is to
51 identify these regions bioinformatically prior to mapping by analysis of sequence complexity (9), or
52 from repetitiveness within the genome (10). Base calls within these pre-specified regions are then
53 ignored ('masking') when assessing relatedness of isolate sequences. A second complementary
54 approach filters base calls based on read mapping confidence as reported by various mappers (11-
55 14) in the form of Mapping Quality (MAQ) scores.

56 *Mycobacterium tuberculosis* is one of the most important pathogens of humans, with about 3 million
57 cases of TB confirmed by culture globally each year (15). Recently, laboratory protocols have been
58 described and deployed by Public Health England (4) in which the species and drug resistance of
59 *Mycobacteria*, including *M. tuberculosis*, are identified by sequencing microbial DNA. Laboratory
60 processing of clinical samples suspected of containing *Mycobacteria* involves decontamination using
61 chemicals which kill non-*Mycobacterial* species before the samples are inoculated into broth culture
62 (16). Mycobacterial Growth Indicator Tubes (MGIT) and associated tube monitoring equipment are
63 a commercially available implementation of such a broth culture system.

64 In the process adopted by Public Health England, sequencing and bioinformatics analysis of DNA
65 extracted from positive MGIT tubes allows determination of Mycobacterial species and drug
66 resistance (17). This laboratory and bioinformatic process also allows the genetic distance between
67 *M. tuberculosis* isolates to be estimated, using sequences derived by consensus base calling from
68 mapped data. The organism co-evolved with human populations as they migrated, generating
69 multiple lineages which differ from the ancestral sequence by hundreds or thousands of single
70 nucleotide polymorphisms (18), as well as small indels, gene deletions and inversions (19).
71 However, the evolutionary clock rate of the organism is slow at about 0.5 single nucleotide variants
72 (SNV)/genome/annum (3, 5, 7) and small numbers of single nucleotide variants (SNV) are of clinical
73 significance: studies based on retrospective collections of *M. tuberculosis*, grown on solid media
74 prior to sequencing, have proposed thresholds of 5 SNVs as compatible with recent transmission (3,
75 5, 7). The bioinformatic processes used for relatedness estimation in the deployed pipeline were

76 also optimised using samples re-grown from frozen stocks on solid media (Lowenstein-Jensen
77 slopes) (16).

78 Quality of complex processes deployed in medical laboratories is assured by adherence to quality
79 standards, such as those laid out in ISO15189 (20). These standards require that the processes
80 followed, and the environments in which they operate, comply with to patterns of work known to
81 enhance the consistency and interpretability of the laboratory outputs. For example, in a drug
82 testing laboratory a set of samples of known composition may be run through the analysers to
83 confirm that particular commonly found substances which might potentially interfere with the assay
84 (such as caffeine, paracetamol, and so on) have no impact on detection of the drug of interest. *M.*
85 *tuberculosis* infection is commonly diagnosed from sputum samples, which contains a wide variety
86 of organisms other than *Mycobacteria* (21). DNA from such organisms may contain sequences
87 homologous with those present in *Mycobacteria*, for example in highly conserved core bacterial
88 genes. As such, this non-*Mycobacterial* DNA has the potential to interfere with assays based on
89 mapping of *Mycobacterial* reference genome mapping.

90 Here, we investigate the concept of interfering substances in the context of the detection of closely
91 related *M. tuberculosis* isolates. In particular, we consider whether DNA of non-*Mycobacterial* origin
92 might cause interference. In doing so, we describe a process which we call *adaptive masking*. This
93 defines 'hard-to-map' regions existing in the context of the laboratory, sequencing and mapping
94 processes being used, independently of predictions based on the reference sequence, and of
95 filtering based on reported mapping quality. Our work was motivated in part by observations from
96 analysis of prospective sequencing of *M. tuberculosis* sequences in England using a previously
97 described bioinformatics pipeline (17). It appeared that high SNV distances were being reported
98 between isolates with a strong epidemiological likelihood of having recently transmitted to each
99 (e.g. isolates with unusual resistance profiles from individuals who were co-habiting): that is, false
100 positive variation between isolates was being reported. We assess the impact of adaptive masking
101 on addressing this problem, and discuss quality control of relatedness monitoring in the context of
102 continuous process monitoring in accredited clinical laboratories.

103

104 METHODS

105 Isolation of DNA from Mycobacteria and Sequencing

106 Clinical specimens were decontaminated and inoculated into Mycobacterial Growth Indicator Tubes
107 (MGIT) tubes (16). Positive samples were extracted (4). After DNA extraction, Illumina sequencing
108 libraries were prepared using either 11 (early in the study) or 15 (later in the study) Mycobacterial
109 DNA extracts, as previously described (4). All samples sent from patients for processing for
110 *Mycobacteria* between 1/05/2016 and 30/05/2017 to a single reference laboratory were studied;
111 the catchment of this laboratory is approximately 15 million people, or about one third of England.

112 Bioinformatic processing

113 Reads obtained from the MiSeq instrument were first examined for the presence of *Mycobacterium*
114 *tuberculosis* using the Mykrobe tool which detects species-specific k-mers (22). Only samples
115 identified as being derived from the *M. tuberculosis* complex by Mykrobe(22) were considered in
116 this work. Additional read classification was performed with Kraken (23), which assigned reads to
117 bespoke database constructed from (i) all bacterial genomes deposited in the NCBI RefSeq database
118 as of January 2017 and (ii) Genome Reference Consortium Human genome build 38 (GRCh38) to
119 allow detection of host DNA as described(24), but with k-mer reduction to 25 gigabases. We
120 quantified reads mapped to *M. tuberculosis* (NCBI taxid: 77643), to non-Mycobacterial bacterial
121 species, and to humans. After doing so, human reads were discarded.

122 Reads were mapped to the H37Rv v2 genome (Genbank NC_000962.2) using Stampy (14), as
123 described (4). Samtools (25) was used to assess sequencing and mapping quality: high quality bases
124 were considered to be those passing the -q30 and -Q30 30 filters (read quality and mapping quality
125 all >30). Consensus sequence was called, requiring a minimum-read depth of 5, including at least
126 one read on each strand. Where an alternative base represented more than 10% of read depth, the
127 base was recorded as uncertain, as described (26).

128 The variant call format (VCF) file was parsed with custom python scripts, and the number of high
129 quality bases (defined using the filters above) at each position determined. These frequencies were
130 extracted, stored and indexed using SQLite using a python API constructed to allow extraction of
131 mixture frequencies in arbitrary positions.

132 Modelling minor variant frequencies

133 We determined the most common (major) variant at each position. All other variants are considered
134 minor variants. We define

135 n - the total sequencing depth at one base;

136 m - the depth of most common variant at one base;
 137 m' - the depth of all variants other than most common, $n-m$ (Suppl. Fig. S1).

138
 139 We divided up the H37Rv genome based on the annotation in NC_000962.3, identifying 8007
 140 regions R, comprising open reading frames and regions between open reading frames
 141 (Supplementary Data S1).

142 For each of these regions $j=1..8007$, if the region has length l_j , the total number of minor variant
 143 bases V_j across the $i=1..l_j$ bases in the region is given by equation (1) and the total read depth D_j
 144 across the gene by equation (2).

$$145 \quad V_j = \sum_{i=1}^{l_j} m'_i \quad (1)$$

$$146 \quad D_j = \sum_{i=1}^{l_j} n_i \quad (2)$$

147
 148 In order to describe the relationship between non-Mycobacterial DNA quantifications and minor
 149 allele frequency, we stratified the number of reads from each sample identified as being from
 150 bacterial genera other than *Mycobacterium* (b) into four approximately equally sized strata:
 151 $b < 1\%$,
 152 $1\% \leq b < 5\%$,
 153 $5\% \leq b < 20\%$,
 154 $b \geq 20\%$.

155
 156 We constructed separate Poisson regression models relating minor base counts (V) for each of the
 157 8,007 regions (with log link and offset $\log(D)$) to the non-Mycobacterial bacterial read categories b
 158 (reference category $< 1\%$), excluding any samples with zero high-quality depth in that region. We
 159 applied Bonferroni correction to model outputs to control for multiple testing ($\alpha = 0.01/8007 =$
 160 1.2×10^{-6}).

161 Comparing impact of masking on pairwise comparisons

162 Based on analysis of model output (see Results), regions with higher minor variant counts than
 163 expected were identified. These regions were excluded from pairwise comparisons performed using
 164 the findNeighbour2 tool (27).

165 Availability of software and data

166 The samples studied have been submitted to NCBI with project number PRJNA451404. Software to
 167 effect the process described, test data, links to the data used, and documentation is freely available
 168 at <http://github.com/davidhwyllie/adaptivemasking>.

169 [Impact of mapper](#)

170 For a random 250 sample subset, we compared the impact of different mappers on minor variant
171 frequencies across regions using the pipeline above. We compared mapping with Stampy
172 1.0.32(14), and Bowtie v 1.2.2 (28) with default parameters, and Bowtie2 v 2.3.4.1 with –very-
173 sensitive and –very-fast preset parameters (29).

174 [Ethical framework](#)

175 Public health action taken as a result of notification and surveillance is one of the Public Health
176 England's key roles as stated in the Health and Social Care Act 2012 and subsequent Government
177 directives which provide the mandate and legislative basis to undertake necessary follow-up. Part of
178 this follow-up is identification of epidemiological and molecular links between cases. This work is
179 part of service development carried out under this framework, and as such explicit ethical approval
180 is unnecessary.

181

182 RESULTS

183 Samples studied

184 The PHE National Mycobacteriology Reference Laboratory Midlands implemented a laboratory
185 process in which specimens received are decontaminated, and inoculated in MGIT bottles, DNA
186 extracts from positive MGIT bottles made, and their contents determined using Illumina short read
187 sequencing (17). Using this process, in the 13 months from 1 May 2016 to 30 May 2017, *M.*
188 *tuberculosis* was identified in 2,751 samples, sent from 2,252 patients (Fig. 1). Using these samples
189 we derived and validated using an independent validation set (Fig.1) a strategy for investigating and
190 controlling false positive variation between samples, which we here term *adaptive masking* (Fig. 2).
191 The initial stages of adaptive masking involve estimating minor variant frequencies across the
192 genome from mapped data, and determining whether these are related to the amount of non-
193 Mycobacterial DNA present.

194 We identified 718 samples from 234 individuals from whom more than one positive sample had
195 been obtained with 7 days of another. Of these, in six individuals samples were either reported as
196 being of different lineages, as defined (30), differed in multiple drug resistances, or differed by >400
197 high quality SNVs. These observations we considered likely due to laboratory or sampling mix-ups,
198 and samples from these patients were excluded. The other 700 samples were used as an
199 independent validation set. From remaining samples, we identified the first sample from each of
200 2,018 individuals which were used to develop the adaptive masking strategy (Fig. 1).

201 Quantifying extraneous DNA and minor variant frequencies post mapping

202 We determined the proportion of non-Mycobacterial bacterial DNA in each sample using Kraken
203 (23), mapped all reads to the H37Rv reference genome irrespective of Kraken results, and filtered
204 the mapped data using stringent quality filters such that the expected error rate is expected to be
205 less than 10^{-3} (see Methods). We defined 8,007 genomic regions in the reference genome; these
206 regions comprise all canonical open reading frames and the genomic regions between them
207 (Supplementary Data S1). We were unable to assess one 15nt region between two PPE family
208 members (positions: 3380439 .. 3380453) as no high quality data mapped there in any sample.

209 In the other 8,006 regions, we observed that both minor variant frequencies and the relationship
210 between minor variant frequency and the number of reads of non-Mycobacterial origin differed
211 markedly by gene. For example, the *B55* and *esxW* genes had respectively very low and very high
212 minor variant frequencies, independent of non-Mycobacterial DNA quantity. A small group of
213 genes, of which the ribosomal component *rrs* is an example, showed low minor variant frequencies,
214 except when non-Mycobacterial DNA was present (Figure 3A, B).

215 Estimating the impact of extraneous bacterial DNA

216 We modelled the relationship between minor variant counts and the number of non-Mycobacterial
217 reads, divided into four approximately equal sized strata (Figure 3B), using Poisson models (Suppl.
218 Data S2). Separate models were constructed for each region. Estimated minor variant frequencies
219 in samples with <1% non-Mycobacterial bacterial reads had a median of 5×10^{-4} (Fig. 4A) across the
220 8006 genomic regions, which is compatible with the expected mapping error rate of $<10^{-3}$, given the
221 filters applied.

222 The distribution of minor variant frequencies when less than 1% non-Mycobacterial bacterial DNA
223 was present approximated a log-normal distribution with mean $\log(5 \times 10^{-4})$ and standard deviation
224 equal to the median absolute deviation (Fig. 4A, observed: black line; fitted, red line), but with a tail
225 to the right. 208 regions (2.6% of the total 8006 regions), including *esxW* as well as other *esx* and
226 PPE family members, had estimated minor variant frequencies $> 2.1 \times 10^{-3}$ when <1% non-
227 Mycobacterial bacterial reads were present (Fig. 4A, Suppl. Data S2). This cutoff represents four
228 median absolute deviations above the median; if the data were log-normally distributed, 24 samples
229 would be expected with minor variant frequencies greater than this, vs. the 208 observed.

230 Overall, estimated minor variant counts rose as non-Mycobacterial DNA concentration rose, but for
231 most regions the increase was small (Fig 4B): the median fold change in minor variant counts in the
232 presence of >20% non-Mycobacterial DNA vs. <1% non-Mycobacterial DNA was 1.097 (i.e. a 9.7%
233 increase, interquartile range 5.4% to 14.0%). 242 regions (3.0%) had statistically significant increases
234 (Fig. 4B, 4C) of more than 50%. Most of these regions had the highest minor variant counts when
235 >20% non-Mycobacterial DNA was present, although a small number had similar minor variant
236 counts in the 5-20% range and the >20% range (Suppl. Fig. S2B).

237 Mutually exclusive regions with increased minor variant frequency

238 Comparing regions with increased minor variant rates with low (<1%) non-Mycobacterial bacterial
239 DNA with those with increased minor variant rates with high (>20%) non-Mycobacterial bacterial
240 DNA shows these regions to be mutually exclusive (Fig. 5). The former include PPE and *esx* family
241 members, while the latter include ribosomal components (*rrl*, *rrs*, *rplB*, and *rps* genes) as well as
242 other highly conserved bacterial genes (tRNA genes, *fusA1*, *infA*, *dnaK*, and others) (Fig. 5, Suppl.
243 Data S2).

244 Examining Kraken read assignments in reads mapped to these highly conserved genes indicated that
245 many reads mapping to these highly conserved regions cannot be unambiguously assigned to the *M.*
246 *tuberculosis* taxon (Supp. Fig. 3) even when non-Mycobacterial bacterial DNA is not present, or is
247 present in small amounts. For example, in the 422 samples for which <1% of *total* bacterial DNA is

248 non-Mycobacterial, only a small proportion (27.3%) of reads mapping with high quality to *rrs* are
249 assigned to the *Mycobacterium* genus by Kraken. However, since very little if any non-Mycobacterial
250 bacterial DNA is present, these reads are almost certainly derived from *M. tuberculosis rrs*. By
251 contrast, the corresponding figure for a gene with little homology with non-Mycobacterial genomes
252 (B55) is 97.4%. The corollary is that if one routinely removes reads which are not assigned by Kraken
253 to the genus of interest (in this case *Mycobacterium*), one will remove a very high proportion of the
254 reads corresponding to critical loci (including drug targets, such as *rrs*) even when no non-
255 Mycobacterial bacterial DNA, as occurs if one was sequencing pure cultures.

256 Adaptive masking reduces the reporting of biologically implausible inter-individual variation

257 A published strategy for excluding regions of high mapping variation within the *M. tuberculosis*
258 genome strategy masks (i.e. excludes from relatedness computations) 277,709 nt (6.3%) of the
259 genome (4). Excluding regions with high estimated minor variant counts with <1% non-
260 Mycobacterial DNA (zone A in Fig. 5B) adds an additional 1.1%. Excluding regions with increased
261 estimated minor variant counts only in the presence of >20% non-Mycobacterial bacterial DNA
262 (zones B-D) masks between 0.2% and 5.1% extra (Fig 5B,C). The masking of regions identified by
263 'adapting' to variation generated during the process forms the final part of the adaptive masking
264 process.

265 In a validation set comprising isolates taken with 7 days of each other from 234 individuals, using the
266 published strategy, 18/346 (5.2%) of pairs studied had ≥ 5 SNV variation, of which 10/346 were \geq
267 20SNV. On exclusion of region D, which comprises the four genes most influenced by non-
268 Mycobacterial DNA, all encoding ribosome associated products (the genes *rrl* and *rrs*, together with
269 the tRNA *metU*, and the highly conserved bacterial gene *tuf*), 0/346 pairs differed by ≥ 5 SNP ($p <$
270 10^{-4} compared with the published method, Wilcoxon test on pairs). These genes were the only
271 genes with minor variant frequencies significantly affected by non-Mycobacterial DNA at the 1 to 5%
272 level (Figure 4B top panel; Supplementary Data Supplementary S2). Additional exclusion of genes in
273 regions B and C, mapping to which is less influenced by non-Mycobacterial DNA, had a limited
274 impact (Figure 6).

275 Impact of mapper

276 In the above work, we use the Stampy mapper(14) which forms part of the deployed TB
277 bioinformatics pipeline(17). To determine whether the choice of mapper was important, we
278 compared the results of mapping a 250 sample test set to the H37Rv reference genome using four
279 different software/parameter settings. We identified regions with any significantly increased
280 mixture frequencies when >20% non-Mycobacterial bacterial DNA was present relative to when <1%

281 was present using a p-value of 0.05, adjusted by Bonferroni's method to 7.6×10^{-6} . The numbers of
282 such regions detected with Stampy, bowtie, and bowtie2 with both stringent and more relaxed
283 matching criteria were 2544, 25, 70, and 72 respectively (Supplementary Data S3). Increased in
284 minor variants frequencies in *rrl* were detected by all techniques; estimated fold increases in minor
285 variant frequencies associated with >20% non-Mycobacterial bacterial DNA were 37, 1.4, 1.5, and
286 1.5 respectively. Thus, although the mismapping observed identified occurs in all cases examined, it
287 is much more prominent with Stampy than with the bowtie series of mappers.

288

289

290

291

292 DISCUSSION

293 Here, we describe an approach which we term adaptive masking. This involves monitoring the
294 minor variant frequency across a bacterial genome to which sequencing reads have been mapped;
295 as such, it measures an end product of next generation sequencing processes, taking into account
296 the natural sequence variation in the samples studied, as well as the impact of DNA extraction and
297 library construction technologies, and the performance of the mapping and filtering software used.

298 Using this approach, we defined a set of 'hard to map' genetic regions (zone A, Fig. 5) with increased
299 minor variant frequencies irrespective of the amount of non-Mycobacterial bacterial DNA.
300 Exclusion of these regions could be considered when assessing consensus *M. tuberculosis* sequences.

301 We also demonstrated a significant positive association between amounts of non-Mycobacterial
302 bacterial DNA and minor variant frequencies in a subset of the mapped genome: significant
303 increases of more than 1.5 fold were observed in 242/8006 regions examined, which together cover
304 about 5% of the *M. tuberculosis* genome. Although the emphasis of this work is on relatedness
305 between isolates, it is notable that included within the 242 regions are a series of genes encoding
306 ribosomal components (*rrs*, *rrl*, *rpoB*, *rpsL*, *rpsA*) which correspond to major antituberculous drug
307 resistance genes(31). Therefore, studies investigating resistance or heteroresistance using these loci
308 should report estimates of the impact of the presence of non-Mycobacterial bacterial DNA on
309 heteroresistance estimates. Such interference may be particularly marked when direct-from-sample
310 short read sequencing is used(32), given the increased ratio of non-Mycobacterial : Mycobacterial
311 DNA in the absence of selective Mycobacterial amplification using culture.

312 Among these 242 regions we identified four 'high variation' regions in which minor variant
313 frequencies are very strongly influenced by non-Mycobacterial bacterial DNA quantities, with fold-
314 increases in minor variant frequencies of >5 in the presence of >20% non-Mycobacterial bacterial
315 DNA. Importantly, if non-Mycobacterial bacterial DNA concentrations are low (<1% of bacterial DNA
316 present), as occurred in retrospective studies when *Mycobacteria* were subcultured on Lowenstein-
317 Jensen slopes prior to sequencing, increased variation is not observed in these regions. The
318 exclusion of the four 'high variation' regions from base calling by a clinically deployed *M.*
319 *tuberculosis* pipeline markedly reduced reported variation between samples derived from the same
320 patient in a short time period. In particular, prior to exclusion of the four high variation regions, in a
321 test set derived from 234 individuals, 5.2% of intra-patient pairs examined differed by 5 SNV or
322 more, with the majority of SNV differences observed in these pairs being >20. Multiple studies
323 indicate this is biologically implausible (3, 5, 7), and after exclusion of the four 'high variation'
324 regions comprising only 0.2% of the genome, no pairs had variation of 5 SNV or more. This suggests

325 that using standard 'masking' and DNA extraction from liquid media false positive variation is
326 reported in a small number of sites in a non-Mycobacterial bacterial DNA dependent manner. Put
327 alternatively, non-Mycobacterial bacterial DNA acts as an interfering substance (20) for relatedness
328 measurements.

329 A potential limitation of this work is that this approach studies pre-specified regions of the genome,
330 specifically coding regions and intergenic regions. This approach was chosen to avoid the challenges
331 of analysing the 4.4×10^6 bases of the *M. tuberculosis* genome individually, with concomitant loss of
332 statistical power. Therefore, as described the method may neither detect, nor allow selective
333 masking of, small regions with high minor variant frequencies within genes. A second limitation is
334 that we did not use of metagenomics classifiers, such as Kraken, to identify non-Mycobacterial
335 'interfering' DNA and eliminate it prior to mapping to the *M. tuberculosis* genome. We did not do
336 this because we observed that for the highly conserved *rrs* genes, metagenomic classifiers cannot
337 confidently assign reads to a genus level, likely because there is insufficient sequence variation
338 within short read sequencing of *rrs* to allow this. Therefore, until longer read sequencing become
339 available sequencing less conserved flanking genomic regions, a strategy of read removal based on
340 metagenomics classification will eliminate a high proportion of *bona fide M. tuberculosis* derived
341 reads in conserved genes, even in samples without any non-Mycobacterial bacterial DNA. Despite
342 these limitations, the strategy chosen appears to be of use clinically, based on the reduction in likely
343 false positive variation between serial samples from individuals.

344 The routine clinical use of next generation sequencing is rapidly increasing (1, 4, 22). However,
345 the reporting of microbial identity, resistotyping, and relatedness information requires complex,
346 multistep processes whose outputs are dependent on specimen decolonisation, selective culture,
347 DNA extraction, library construction, DNA sequencing, and bioinformatic analysis (4). Reagent
348 batches, software versions, and equipment involved in the process are all subject to change over
349 time. The adaptive masking approach we describe here represents a route to the quantitative
350 monitoring the performance of the output of this pathway, identifying whether changes in process
351 which may appear innocuous alter mapping and basecalling across the genome. We do not
352 propose that the output from the adaptive masking process as demonstrated here with data
353 generated by Public Health England and processed by particular bioinformatics tools should be used
354 to generate a list of problematic genomic regions which can be universally applied. Rather, we
355 envisage that the adaptive masking process will be performed as part of the acceptance of process
356 change, and periodically as part of quality monitoring, under the exact conditions used in the clinical
357 laboratory issuing NGS based results. A list of problematic positions to be ignored during

358 relatedness calculations can then be fed into systems doing such calculations, such as
359 findNeighbour2 (27), which would apply such masking across all samples.

360 Generalizable to other organisms and mapping pipelines, the adaptive masking approach we
361 described here will have application in monitoring the performance of such processes quantitatively,
362 in interpreting estimates of possible heteroresistance, and in preventing the calling of false positive
363 variation in the context of clinically deployed genomics.

364

365 ACKNOWLEDGEMENTS

366 This study is supported by the Health Innovation Challenge Fund (a parallel funding partnership
367 between the Wellcome Trust [WT098615/Z/12/Z] and the Department of Health [grant HICF-T5-
368 358]) and NIHR Oxford Biomedical Research Centre. Professor Derrick Crook is affiliated to the
369 National Institute for Health Research Health Protection Research Unit (NIHR HPRU) in Healthcare
370 Associated Infections and Antimicrobial Resistance at University of Oxford in partnership with Public
371 Health England. Professor Crook is based at University of Oxford. The views expressed are those of
372 the author(s) and not necessarily those of the NHS, the NIHR, the Department of Health or Public
373 Health England. The sponsors of the study had no role in study design, data collection, data analysis,
374 data interpretation, or writing of the report. The corresponding author had full access to all the data
375 in the study and had final responsibility for the decision to submit for publication.

376

377

378 REFERENCES

- 379 1. De Silva D, Peters J, Cole K, Cole MJ, Cresswell F, Dean G, Dave J, Thomas DR, Foster K,
380 Waldram A, Wilson DJ, Didelot X, Grad YH, Crook DW, Peto TE, Walker AS, Paul J, Eyre DW.
381 2016. Whole-genome sequencing to determine transmission of *Neisseria gonorrhoeae*: an
382 observational study. *Lancet Infect Dis* 16:1295-1303.
- 383 2. Gordon NC, Pichon B, Golubchik T, Wilson DJ, Paul J, Blanc DS, Cole K, Collins J, Cortes N,
384 Cubbon M, Gould FK, Jenks PJ, Llewelyn M, Nash JQ, Orendi JM, Paranthaman K, Price JR,
385 Senn L, Thomas HL, Wyllie S, Crook DW, Peto TEA, Walker AS, Kearns AM. 2017. Whole-
386 Genome Sequencing Reveals the Contribution of Long-Term Carriers in *Staphylococcus*
387 *aureus* Outbreak Investigation. *J Clin Microbiol* 55:2188-2197.
- 388 3. Walker TM, Lalor MK, Broda A, Ortega LS, Morgan M, Parker L, Churchill S, Bennett K,
389 Golubchik T, Giess AP, Del Ojo Elias C, Jeffery KJ, Bowler I, Laurenson IF, Barrett A,
390 Drobniewski F, McCarthy ND, Anderson LF, Abubakar I, Thomas HL, Monk P, Smith EG,
391 Walker AS, Crook DW, Peto TEA, Conlon CP. 2014. Assessment of *Mycobacterium*
392 *tuberculosis* transmission in Oxfordshire, UK, 2007-12, with whole pathogen genome
393 sequences: an observational study. *Lancet Respir Med* 2:285-292.
- 394 4. Pankhurst LJ, Del Ojo Elias C, Votintseva AA, Walker TM, Cole K, Davies J, Fermont JM,
395 Gascoyne-Binzi DM, Kohl TA, Kong C, Lemaitre N, Niemann S, Paul J, Rogers TR, Roycroft E,
396 Smith EG, Supply P, Tang P, Wilcox MH, Wordsworth S, Wyllie D, Xu L, Crook DW. 2016.
397 Rapid, comprehensive, and affordable mycobacterial diagnosis with whole-genome
398 sequencing: a prospective study. *Lancet Respir Med* 4:49-58.
- 399 5. Bryant JM, Harris SR, Parkhill J, Dawson R, Diacon AH, van Helden P, Pym A, Mahayiddin AA,
400 Chuchottaworn C, Sanne IM, Louw C, Boeree MJ, Hoelscher M, McHugh TD, Bateson AL,
401 Hunt RD, Mwaigwisya S, Wright L, Gillespie SH, Bentley SD. 2013. Whole-genome sequencing
402 to establish relapse or re-infection with *Mycobacterium tuberculosis*: a retrospective
403 observational study. *Lancet Respir Med* 1:786-92.
- 404 6. Casali N, Nikolayevskyy V, Balabanova Y, Harris SR, Ignatyeva O, Kontsevaya I, Corander J,
405 Bryant J, Parkhill J, Nejentsev S, Horstmann RD, Brown T, Drobniewski F. 2014. Evolution and
406 transmission of drug-resistant tuberculosis in a Russian population. *Nat Genet* 46:279-86.
- 407 7. Guerra-Assuncao JA, Crampin AC, Houben RM, Mzembe T, Mallard K, Coll F, Khan P, Banda L,
408 Chiwaya A, Pereira RP, McNerney R, Fine PE, Parkhill J, Clark TG, Glynn JR. 2015. Large-scale
409 whole genome sequencing of *M. tuberculosis* provides insights into transmission in a high
410 prevalence area. *Elife* 4.
- 411 8. Guerra-Assuncao JA, Houben RM, Crampin AC, Mzembe T, Mallard K, Coll F, Khan P, Banda L,
412 Chiwaya A, Pereira RP, McNerney R, Harris D, Parkhill J, Clark TG, Glynn JR. 2015. Recurrence
413 due to relapse or reinfection with *Mycobacterium tuberculosis*: a whole-genome sequencing
414 approach in a large, population-based cohort with a high HIV infection prevalence and active
415 follow-up. *J Infect Dis* 211:1154-63.
- 416 9. Morgulis A, Gertz EM, Schaffer AA, Agarwala R. 2006. A fast and symmetric DUST
417 implementation to mask low-complexity DNA sequences. *J Comput Biol* 13:1028-40.
- 418 10. Bedell JA, Korf I, Gish W. 2000. MaskerAid: a performance enhancement to RepeatMasker.
419 *Bioinformatics* 16:1040-1.
- 420 11. Langmead B, Trapnell C, Pop M, Salzberg SL. 2009. Ultrafast and memory-efficient alignment
421 of short DNA sequences to the human genome. *Genome Biol* 10:R25.
- 422 12. Li H, Ruan J, Durbin R. 2008. Mapping short DNA sequencing reads and calling variants using
423 mapping quality scores. *Genome Res* 18:1851-8.
- 424 13. Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler
425 transform. *Bioinformatics* 25:1754-60.

- 426 14. Lunter G, Goodson M. 2011. Stampy: a statistical algorithm for sensitive and fast mapping of
427 Illumina sequence reads. *Genome Res* 21:936-9.
- 428 15. Vijay S, Dalela G. 2016. Prevalence of LRTI in Patients Presenting with Productive Cough and
429 Their Antibiotic Resistance Pattern. *J Clin Diagn Res* 10:Dc09-12.
- 430 16. Anonymous. 2017. SMI B 40: Investigation of specimens for Mycobacterium species. Public
431 Health England, London.
- 432 17. Quan TP, Bawa Z, Foster D, Walker T, Del Ojo Elias C, Rathod P, Iqbal Z, Bradley P, Mowbray
433 J, Walker AS, Crook DW, Wyllie DH, Peto TEA, Smith EG. 2017. Evaluation of whole genome
434 sequencing for Mycobacterial species identification and drug susceptibility testing in a
435 clinical setting: a large-scale prospective assessment of performance against line-probe
436 assays and phenotyping. *J Clin Microbiol* doi:10.1128/jcm.01480-17.
- 437 18. Comas I, Coscolla M, Luo T, Borrell S, Holt KE, Kato-Maeda M, Parkhill J, Malla B, Berg S,
438 Thwaites G, Yeboah-Manu D, Bothamley G, Mei J, Wei L, Bentley S, Harris SR, Niemann S,
439 Diel R, Aseffa A, Gao Q, Young D, Gagneux S. 2013. Out-of-Africa migration and Neolithic
440 coexpansion of Mycobacterium tuberculosis with modern humans. *Nat Genet* 45:1176-82.
- 441 19. Hatherell HA, Colijn C, Stagg HR, Jackson C, Winter JR, Abubakar I. 2016. Interpreting whole
442 genome sequencing for investigating tuberculosis transmission: a systematic review. *BMC*
443 *Med* 14:21.
- 444 20. Anonymous. 2012. ISO15189:2012 Medical Laboratories - Requirements for Quality and
445 Competence. International Organisation for Standardisation.
- 446 21. Nasidze I, Li J, Quinque D, Tang K, Stoneking M. 2009. Global diversity in the human salivary
447 microbiome. *Genome Research* 19:636-643.
- 448 22. Bradley P, Gordon NC, Walker TM, Dunn L, Heys S, Huang B, Earle S, Pankhurst LJ, Anson L,
449 de Cesare M, Piazza P, Votintseva AA, Golubchik T, Wilson DJ, Wyllie DH, Diel R, Niemann S,
450 Feuerriegel S, Kohl TA, Ismail N, Omar SV, Smith EG, Buck D, McVean G, Walker AS, Peto TE,
451 Crook DW, Iqbal Z. 2015. Rapid antibiotic-resistance predictions from genome sequence
452 data for Staphylococcus aureus and Mycobacterium tuberculosis. *Nat Commun* 6:10063.
- 453 23. Wood DE, Salzberg SL. 2014. Kraken: ultrafast metagenomic sequence classification using
454 exact alignments. *Genome Biol* 15:R46.
- 455 24. Street TL, Sanderson ND, Atkins BL, Brent AJ, Cole K, Foster D, McNally MA, Oakley S, Peto L,
456 Taylor A, Peto TEA, Crook DW, Eyre DW. 2017. Molecular Diagnosis of Orthopedic-Device-
457 Related Infection Directly from Sonication Fluid by Metagenomic Sequencing. *J Clin*
458 *Microbiol* 55:2334-2347.
- 459 25. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R.
460 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25:2078-9.
- 461 26. Walker TM, Kohl TA, Omar SV, Hedge J, Del Ojo Elias C, Bradley P, Iqbal Z, Feuerriegel S,
462 Niehaus KE, Wilson DJ, Clifton DA, Kapatai G, Ip CL, Bowden R, Drobniewski FA, Allix-Beguec
463 C, Gaudin C, Parkhill J, Diel R, Supply P, Crook DW, Smith EG, Walker AS, Ismail N, Niemann S,
464 Peto TE. 2015. Whole-genome sequencing for prediction of Mycobacterium tuberculosis
465 drug susceptibility and resistance: a retrospective cohort study. *Lancet Infect Dis* 15:1193-
466 202.
- 467 27. Mazariegos-Canellas O, Do T, Peto T, Eyre DW, Underwood A, Crook D, Wyllie DH. 2017.
468 BugMat and FindNeighbour: command line and server applications for investigating bacterial
469 relatedness. *BMC Bioinformatics* 18:477.
- 470 28. Langmead B, Trapnell C, Pop M, Salzberg SL. 2009. Ultrafast and memory-efficient alignment
471 of short DNA sequences to the human genome. *Genome Biology* 10:R25.
- 472 29. Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nature Methods*
473 9:357.
- 474 30. Coll F, McNerney R, Guerra-Assuncao JA, Glynn JR, Perdigao J, Viveiros M, Portugal I, Pain A,
475 Martin N, Clark TG. 2014. A robust SNP barcode for typing Mycobacterium tuberculosis
476 complex strains. *Nat Commun* 5:4812.

- 477 31. Walker TM, Kohl TA, Omar SV, Hedge J, Del Ojo Elias C, Bradley P, Iqbal Z, Feuerriegel S,
478 Niehaus KE, Wilson DJ, Clifton DA, Kapatai G, Ip CLC, Bowden R, Drobniowski FA, Allix-Béguec
479 C, Gaudin C, Parkhill J, Diel R, Supply P, Crook DW, Smith EG, Walker AS, Ismail N, Niemann S,
480 Peto TEA. Whole-genome sequencing for prediction of *Mycobacterium*
481 tuberculosis drug susceptibility and resistance: a retrospective cohort study. The
482 Lancet Infectious Diseases 15:1193-1202.
- 483 32. Votintseva AA, Bradley P, Pankhurst L, Del Ojo Elias C, Loose M, Nilgiriwala K, Chatterjee A,
484 Smith EG, Sanderson N, Walker TM, Morgan MR, Wyllie DH, Walker AS, Peto TEA, Crook DW,
485 Iqbal Z. 2017. Same-Day Diagnostic and Surveillance Data for Tuberculosis via Whole-
486 Genome Sequencing of Direct Respiratory Samples. J Clin Microbiol 55:1285-1298.
- 487
- 488

489 **FIGURE LEGENDS**

490 **Figure 1 Samples used, derivation and validation sets**

491 A flow chart describing the samples used, and the selection of derivation and validation sets.

492 **Figure 2 Bioinformatic Processes**

493 A flow diagram illustrating the standard bioinformatics pipeline used, as well as the adaptive
494 masking process used to generate masks. Gray circles denote links to a description of the process at
495 <https://github.com/davidhwyllie/adaptivemasking>. FindNeighbour2 is an open-source server-based
496 system for monitoring single nucleotide variation(27).

497 **Figure 3 Minor variant frequency and non-Mycobacterial bacterial DNA quantities**

498 The observed minor variant frequency for three regions of the *M. tuberculosis* genome (genes *B55*,
499 *eswX*, and *rrs*) vs. the proportion of reads of non-Mycobacterial bacterial origin (as determined by
500 Kraken) is shown for samples in the derivation set (n=2018). Panel (A) shows a dot plot, whereas in
501 panel (B) the proportion of reads of non-Mycobacterial bacterial origin is stratified into strata with
502 1%, 5% and 20% boundaries. Numbers above each stratum refer to the number of samples with
503 non-zero read depth in that region.

504 **Figure 4 Modelling minor variant frequencies**

505 For 8,006 genomic regions of the H37Rv reference genome, Poisson models were used to estimate
506 the mean minor variant frequency. The estimated minor variant frequency when less than 1% non-
507 Mycobacterial bacterial DNA is present (n=208, 2.6% regions) is shown in (A). The red line is a log
508 normal distribution with $\mu = \log(\text{minor variant frequency with } <1\% \text{ non-Mycobacterial DNA})$ and $\sigma =$
509 $\text{median absolute deviation}(\log(\text{minor variant frequency with } <1\% \text{ non-Mycobacterial DNA}))$. In (B)
510 the rate ratio estimates (i.e. the fold change associated with increases in non-Mycobacterial
511 bacterial DNA quantifications) for each gene are shown. (C) shows the significance of a test
512 comparing the $\log(\text{rate ratio estimates})$ with zero, in the form of a Volcano plot. The dotted lines in
513 (B) and (C) correspond to a 50% increase in rate ratio.

514

515 **Figure 5 A distinct subset of genes are impacted by quantity of non-Mycobacterial DNA**

516 Legend: (A) Fold change in minor variant frequency with > 20% non-Mycobacterial bacterial DNA
517 present vs. <1% non-Mycobacterial bacterial DNA. Quadrant boundary markers correspond to
518 (horizontal line) a 50% increase over <1% non-Mycobacterial bacterial DNA and (vertical line) a
519 minor variant frequency of 2.1×10^{-3} . (B) Genes with elevated minor variant frequencies when non-
520 Mycobacterial bacterial DNA is low (<1%) or high (>20%) fall into mutually exclusive sets. (C) the

521 number of bases represented by the deployed masking, vs. the deployed masking plus the genes in
522 zones A, A+D, A+D+C, and A+D+C+B.

523 [Figure 6 Impact of masking strategies on reported distances between closely related samples](#)
524 SNV distances between pairs of *M. tuberculosis* genomes isolates from samples taken from the same
525 individual within 7 days of each other were compared using different masking strategies. The top
526 panel describes the published, deployed method of masking. In the panels below, genes in the
527 zones shown in Figure 5B are additionally masked (i.e. ignored from pairwise comparisons).

528











